

# LLM Alignment Fine-tuning via Human Preference: DPO vs IPO

Tri Nguyen

nguyetr9@oregonstate.edu

March 27, 2025

Large language model (LLM) alignment is a crucial step in increasing LLM’s usability and safety. This topic has been received a lot of attentions, resulting in very interesting development over the last 5 years. I’m particularly interested in a line of development including 3 works:

- **Reinforcement Learning from Human Feedback (RLHF)** [5, 4, 2, 8] from OpenAI: learning score model, using RL to perform alignment
- **Direct preference optimization (DPO)** [6] from Stanford: learning score model, alignment is “automatically” obtained.
- **IPO (or  $\Psi$ PO )** [1] from Google DeepMind: no need for learning score model, at all :))

is a very nice buildup as one addresses the previous work’s issues reducing learning process’ complexity. An overview comparing the three method is shown in Figure 1.

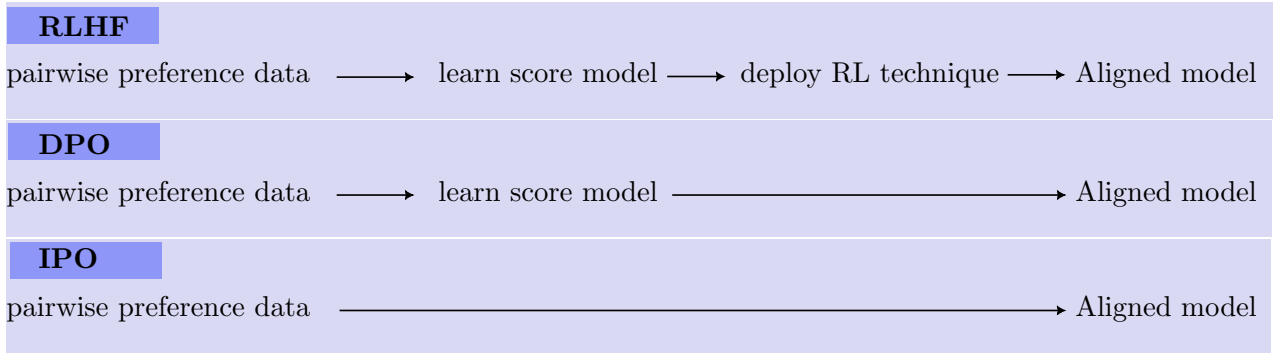


Figure 1: An overview comparison of RLHF, DPO, and IPO.

In the first 2 works, both RLHF and DPO’s idea are elegant and simple to grasp. This is however not the case for IPO. When I first read the IPO paper the whole thing screams “ad-hoc” and *unnecessarily-complicated* ;). Their derivation is a bit confusing, their empirical loss looks counter-intuitive, they don’t even have any LLM alignment experiments. But the fact that IPO method works well empirically (from my own experience as well as from other papers) bothers me quite a lot. After spending some effort to examining this method more carefully, IPO’s idea turns out to be quite nice and very clever. If you don’t get it from skimming through the paper (like I did), I hope this blog post can convince you to have another look at this method.

## 1 LLM alignment

LLMs obtained after a pre-training step over vast un-labeled datasets possesses an amazingly ability of natural-looking text completions. However, due to the nature of unsupervised training, the resulting

models exhibit many undesired behaviors (via the generation), including being unhelpful, biased, sexist, racism, hallucination.

LLM alignment aims to address this issue by steering model’s behavior toward the desired characterizations using following formulation

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{y} \sim \pi} [s(\mathbf{x}, \mathbf{y}) - \beta D_{\text{kl}}(\pi \parallel \pi_{\text{ref}})] \right] \quad (1)$$

Here, the so-called *score function*  $s(\mathbf{x}, \mathbf{y})$  is assumed to produce a scalar value indicating how strongly the response  $\mathbf{y}$  is aligned with the desired characteristic given the prompt  $\mathbf{x}$ . We wish find a policy (a language model)  $\pi$  such that it retains most of the good text-generation capability of  $\pi_{\text{ref}}$ , and at the same time producing responses  $\mathbf{y}$  to maximize the score function. The balance between 2 objectives is controlled by  $\beta$ . A too small  $\beta$  might lead to overly optimized  $\pi$  that loses the generally good text generation of  $\pi_{\text{ref}}$ , while too large  $\beta$  prevents  $\pi$  from adjusting toward better alignment. The reference policy  $\pi_{\text{ref}}$  is can be seen as a good initialization.

An obvious barrier in solving (1) is the unknown  $s(\mathbf{x}, \mathbf{y})$ . It is very non-trivial to hand-crafting the score function  $s(\mathbf{x}, \mathbf{y})$ . Instead, a popular approach is to learn  $s(\mathbf{x}, \mathbf{y})$  using pairwise preference datasets. A sample of such dataset consists of a tuple  $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2, c)$  where  $\mathbf{x}$  is a prompt,  $\mathbf{y}_1, \mathbf{y}_2$  are 2 possible responses (a continuation of prompt  $\mathbf{x}$ ). These three elements are sent to a human annotator, who assigns a label  $c \in \{1, 2\}$  to indicate which response is more preferred with the alignment objective. For instance,  $c = 1$  implies that  $\mathbf{y}_1$  is preferred over  $\mathbf{y}_2$ , denoting as  $\mathbf{y}_1 \succ \mathbf{y}_2$ .

Different methods propose to solve (1) in different ways but they all share the same principle: using the preference data to “infer” the unknown function  $s(\mathbf{x}, \mathbf{y})$ . As the only available supervised signal is the preference dataset, it is necessary to assume certain specification to relate the unknown score function  $s(\mathbf{x}, \mathbf{y})$  and the collected pairwise preference data.

$$\text{Unknown score function } s(\mathbf{x}, \mathbf{y}) \xleftrightarrow{\text{certain specification}} \text{Preference data } (\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2, c) \quad (2)$$

In the following, we will go through different realizations of this relations, leading to 3 different popular techniques: RLHF, DPO, and IPO.

## 2 RLHF and DPO

### 2.1 RLHF

The structure of the optimization problem (1) is very much a typical reinforcement learning (RL) problem, except that the score function  $s(\mathbf{x}, \mathbf{y})$  is unknown. Naturally, if one can learn  $s(\mathbf{x}, \mathbf{y})$ , then an off-the-shelf RL method can be deployed to solve (1). This is the very idea proposed by RLHF.

To learn  $s(\mathbf{x}, \mathbf{y})$ , RLHF assumes the Bradley-Terry (BT) model [3] for the preference data generation, which hypothesizes that

$$\Pr(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x}) \triangleq \Pr(c = 1 \mid \mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = \frac{\exp(s(\mathbf{x}, \mathbf{y}_1))}{\exp(s(\mathbf{x}, \mathbf{y}_1)) + \exp(s(\mathbf{x}, \mathbf{y}_2))} = \sigma(s(\mathbf{x}, \mathbf{y}_1) - s(\mathbf{x}, \mathbf{y}_2))$$

Since the score function gives a higher value for a better aligned response, it is more likely that that response is preferred over the other.

We can see that this assumption specifies the relation as mentioned in (2). If we further assume the true  $s(\mathbf{x}, \mathbf{y})$  belong to certain class of deep neural network (such as an LLM-based network), then the true score function can be recovered using MLE:

$$\underset{\theta}{\text{minimize}} \quad - \mathbb{E}_{(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2), c \sim \mathcal{D}} [\mathcal{L}_{\text{logistic}}(\sigma(s_{\theta}(\mathbf{x}, \mathbf{y}_1) - s_{\theta}(\mathbf{x}, \mathbf{y}_2)), c)] \quad (3)$$

where

$$\mathcal{L}_{\text{logistic}}(p, c) = \mathbb{I}[c = 1] \log p + \mathbb{I}[c = 2] \log (1 - p).$$

In practice,  $s_{\theta}(\mathbf{x}, \mathbf{y})$  can be parameterized using another large language model. After obtaining the optimal solution  $\theta^*$  to (3), the alignment fine-tuning is performed by plugging  $s_{\theta^*}$  into (1) and an off-the-shelf RL techniques, such as PPO [7] is invoked to solve (1).

This approach while being straightforward, suffers from several technical difficulties:

- A 2-stage training pipeline is complex and prone to error accumulation.
- The use of RL require intensive and careful hyperparameter tuning.

These challenges are the main motivations for the development of DPO.

## 2.2 DPO

DPO improves upon RLHF by eliminating the RL step. In particular, DPO’s authors realizes that under the same preference model (BT model), the relation between preference label and score function only depends on the relative difference in score, not the absolute score values. This enables them to use a clever trick to re-parameterize the score function via an optimal policy, effectively eliminate the needs of deploying RL.

Notice that the objective in (1) can be expressed as:

$$\mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [s(\mathbf{y}, \mathbf{x}) - \beta D_{\text{kl}}(\pi \parallel \pi_{\text{ref}})] \right] = \mathbb{E}_{\mathbf{x}} \left[ D_{\text{kl}}(\pi \parallel \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp(\beta^{-1} s(\mathbf{y}, \mathbf{x}))) \right] + \text{const}$$

where  $Z(\mathbf{x})$  is an intractable normalizing factor. As KL-divergence reaches minimum value at 0, this expression suggests an optimal solution  $\pi^*$  as

$$\pi^*(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp(\beta^{-1} s(\mathbf{y}, \mathbf{x})),$$

which equivalently implies

$$s(\mathbf{y}, \mathbf{x}) = \beta \log \frac{\pi^*(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} + \beta \log Z(\mathbf{x})$$

This identity establishes a relation between an arbitrary score function  $s(\mathbf{x}, \mathbf{y})$  and a corresponding optimal policy  $\pi^*(\mathbf{y} | \mathbf{x})$  with respect to that score function. It is not very useful by itself due to the intractable factor  $Z(\mathbf{x})$ . However, the relative score difference, which is all that matters, is **independent** of  $Z(\mathbf{x})$ :

$$s(\mathbf{x}, \mathbf{y}_1) - s(\mathbf{x}, \mathbf{y}_2) = \beta \log \frac{\pi^*(\mathbf{y}_1 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1 | \mathbf{x})} - \beta \log \frac{\pi^*(\mathbf{y}_2 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2 | \mathbf{x})}. \quad (4)$$

To be more concise, define

$$h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = \beta \log \frac{\pi_{\theta}(\mathbf{y}_1 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1 | \mathbf{x})} - \beta \log \frac{\pi_{\theta}(\mathbf{y}_2 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2 | \mathbf{x})}.$$

then equation (4) can be shorten as

$$s(\mathbf{x}, \mathbf{y}_1) - s(\mathbf{x}, \mathbf{y}_2) = h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2). \quad (5)$$

This condition plays a key in ensuring policy  $\pi_{\theta}$  being the optimal solution to the original alignment formulation (1). Particularly, it is shown in DPO that any policy  $\pi_{\theta}$  satisfying condition (5) is an optimal solution to the alignment formula in (1), up to some trivial ambiguity.

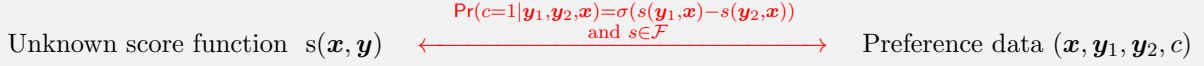
With this insight, DPO proposed to use  $h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$  to parameterize the relative score difference between 2 responses  $\mathbf{y}_1, \mathbf{y}_2$ . As before, they employ the BT model and use MLE to derive the loss function:

$$\underset{\theta}{\text{minimize}} \quad - \mathbb{E}_{(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2), c \sim \mathcal{D}} [\mathcal{L}_{\text{logistic}}(\sigma(h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)), c)]. \quad (6)$$

With this parameterization, an optimal solution  $\theta^*$  to (6) gives us an optimal policy  $\pi_{\theta^*}$  to (1) simultaneously.

## 2.3 The Common Theme

In terms of modeling, both RLHF and DPO relies on the same specification:



There are 2 factors in this specification:

- The BT model is used to relate the score function and the preference data
- The score function is assumed to belong to certain known hypothesis class  $\mathcal{F}$ , either an arbitrary neural network as in RLHF, or a structured neural network as in DPO.

The combination of the two enables learning  $s(\mathbf{x}, \mathbf{y})$  using preference data.

### Drawback.

- The preference data relies on BT model. BT model is intuitive and is successfully deployed in many domains, such as economy, however, it is still restrictive and who knows if real data generation truly follows it.
- The particular functional form of BT model, i.e., the sigmoid function makes it computationally difficult to model deterministic cases. Specifically, to have  $\Pr(c = 1 | \mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) \rightarrow 1$ , the quantity  $s(\mathbf{x}, \mathbf{y}_1) - s(\mathbf{x}, \mathbf{y}_2)$  in RLHF, or  $h_\theta(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$  in DPO needs to reach  $\infty$ . This behavior causes a particular detrimental effect: worsening the reward hacking in the second phase in RLHF, or overfitting in DPO where the learned policy drifting arbitrarily far away from  $\pi_{\text{ref}}$  regardless of  $\beta$ .
- And subtly, the alignment fine-tuning task is cast into a score learning problem. This is in turn solved indirectly. What we wish to estimate is in the “logits” space, i.e.,  $\hat{h}(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \approx s(\mathbf{y}_1, \mathbf{x}) - s(\mathbf{y}_2, \mathbf{x})$  but what the criterion enforcing is in the probability space, i.e.,  $\sigma(\hat{h}(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})) \approx \sigma(s(\mathbf{y}_1, \mathbf{x}) - s(\mathbf{y}_2, \mathbf{x}))$ . A small mismatch in the later could result in a much large mismatch in the former estimation, which affects the ultimate alignment task. This is, to me, the major instability issue of RLHF and DPO.

And these drawback lead to the development of IPO.

## 3 IPO

The method in RLHF/DPO can be seen to boil down to a binary classification problem: Given a sample  $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$ , one wish to predict if the label is 1 or 2. Their loss functions are very intuitive following this view. Take DPO for example, if annotator is very certain that the label  $c = 1 \Leftrightarrow \mathbf{y}_1 \succ \mathbf{y}_2$ , the score gap  $h_\theta(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})$  between the 2 responses should be large. The loss functions DPO (6) promotes this by increasing the estimated likelihood  $\sigma(h_\theta(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}))$ . Similar mechanism can be observed in RLHF as well.

This very intuitive idea might have spoiled me so much that when I first encountered the IPO paper, every step of it felt wrong. Let me articulate. Let’s look at their final loss:

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell} \left[ (h_\theta(\mathbf{y}_w, \mathbf{y}_\ell, \mathbf{x}) - 0.5)^2 \right],$$

where notation  $\mathbf{y}_w, \mathbf{y}_\ell$  is an equivalent form of using the label  $c \in \{1, 2\}$  and imply  $\mathbf{y}_w \succ \mathbf{y}_\ell$ . This is nothing special and everyone uses it.

Now come to the unintuitive parts:

- Why regression? How can it turn a classification to a regression problem?

- And why regressing toward 0.5?
- Furthermore, IPO claims that they even **don't need to learn the score function**. Then how do they do alignment at all!!!

These baffling points have been in my mind for a while. Unfortunately, their derivation in the paper didn't help much. Recently, I have some troubles comparing against this IPO method, and thus I've spent some more time on it. Turned out, their solution is more elegant than I thought.

### 3.1 IPO's Loss Derivation

Note that IPO's goal is to perform the alignment fine-tuning task via the same formulation (1) as RLHF and DPO do. Unlike previous methods where there is no structure on  $s(\mathbf{x}, \mathbf{y})$ , IPO assumes the following (unknown) score function:

$$s^{\dagger}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{y}' \sim \mu(\cdot | \mathbf{x})} [v(\mathbf{y}, \mathbf{y}', \mathbf{x})]. \quad (7)$$

The scalar-valued function  $v(\mathbf{y}, \mathbf{y}', \mathbf{x}) \in [0, 1]$  denotes the probability of  $\mathbf{y} \succ \mathbf{y}'$ . This score function can be seen as a quantification of how a response  $\mathbf{y}$  is preferred on average with respect to a predefined policy  $\mu$ . For now,  $\mu$  is just an arbitrary policy. Using this score function, the alignment task is to find a policy  $\pi^*$  within the proximity of  $\pi_{\text{ref}}$  and be better than policy  $\mu$  as much as possible. In essence, the goal is quite the same, except that we have certain structure on the score function to work with. As we will see, using this unknown score function, IPO cleverly perform alignment without explicitly learning it.

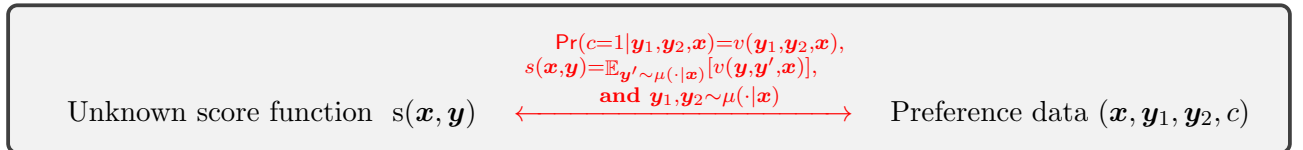
Note that in the paper, IPO denote  $v(\mathbf{y}, \mathbf{y}', \mathbf{x}) \triangleq p(\mathbf{y} \succ \mathbf{y}' | \mathbf{x})$ . Although using  $p(\mathbf{y} \succ \mathbf{y}' | \mathbf{x})$  can be intuitive, I feel much comfortable using  $v(\mathbf{y}, \mathbf{y}', \mathbf{x})$  to explicitly remind myself that it is nothing but a special unknown function. As  $v$  denote a probability of preference, we require it to satisfy:

$$v(\mathbf{y}, \mathbf{y}', \mathbf{x}) \in [0, 1], \quad \forall \mathbf{y}, \mathbf{y}', \mathbf{x} \quad (8a)$$

$$v(\mathbf{y}, \mathbf{y}, \mathbf{x}) = 0.5, \quad \forall \mathbf{y}, \mathbf{x} \quad (8b)$$

$$v(\mathbf{y}, \mathbf{y}', \mathbf{x}) = 1 - v(\mathbf{y}', \mathbf{y}, \mathbf{x}) \quad (8c)$$

Condition (8a) is to ensure valid probabilities, condition (8b) and (8c) are to enforce the physical meaning of pairwise preference. With these notations, IPO's specification on the relation between preference labels and score function is as follows



Notice that  $v$  is still unknown but IPO don't aim to learn  $v$ . The new, and important requirement is that the pair of responses are assumed to be drawn from the same behavior policy  $\mu$ . Recall that  $\mu$  plays a role in defining the score function. This seems to be a limitation. However, it is possible that particular choice of  $\mu$  does not change the optimal policy much (or at all).

As in DPO, IPO enforces the condition (5) to ensure optimal policy:

$$h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = s(\mathbf{x}, \mathbf{y}_1) - s(\mathbf{x}, \mathbf{y}_2) = \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})],$$

which can be realized by solve an optimization problem:

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[ \left( h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) - \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right)^2 \right] \quad (9)$$

Notice that the regression target is unknown. However, by exploiting 2 key factors in their specification: (i) the pair of responses are drawn from the same policy  $\mu$  which is also used in defining the score

function, and (ii) the structure of  $h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$ , IPO shows that optimization problem (9) is equivalent to

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[ (h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) - v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}))^2 \right]. \quad (10)$$

I have a more rigorous derivation in the pdf attachment proving this part (Section ??). The derivation is not super complex but it certainly is not trivial. I am very baffled as how they came up with this observation.

The equivalent optimization (10) is tremendously useful since the target in (10) can be approximated as

$$v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \approx c - 1$$

With such approximation, IPO proposed the following optimization problem:

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mu, c} \left[ (h_{\theta}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) - c + 1)^2 \right]. \quad (11)$$

This can be further simplified when considering symmetrical property of the preference data. For instance, a sample  $(\mathbf{y}, \mathbf{y}', \mathbf{x}, 1)$  induces another sample  $(\mathbf{y}', \mathbf{y}, \mathbf{x}, 2)$ . Exploiting this structure and denote  $\mathbf{y}_w, \mathbf{y}_\ell$  based on the label  $c$  such that  $\mathbf{y}_w \succ \mathbf{y}_\ell$ , the previous problem is equivalent to

$$\underset{\theta}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell \sim \mu} \left[ (h_{\theta}(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_\ell) - 0.5)^2 \right]. \quad (12)$$

### 3.2 Take-Home Points

To answer the questions in the beginning on this section:

- How can it turn a classification to a regression problem? It does not. RLHF/DPO learn the underlying model by casting the problem as a classification problem. IPO does not learn any underlying model, and hence there is no need of any classification problem. The regression problem stems from enforcing condition (5) for the optimal policy.
- Why regressing toward 0.5? This fixed number 0.5 is a result to a quite rough approximation  $\Pr(c = 1 \mid \mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \approx c - 1$ .
- Lastly, as IPO's working directly on the logit space, it seems to be more numerically robust compared to previous methods.

### 3.3 Technical Derivation

**Lemma 3.1.** *Optimization problems (9) and (10) are equivalent.*

*Proof.* Define

$$\mathbf{q}_{\theta}(\mathbf{y}, \mathbf{x}) = \beta \log \frac{\pi_{\theta}(\mathbf{y} \mid \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} \mid \mathbf{x})},$$

then  $h_{\theta}(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) = \mathbf{q}_{\theta}(\mathbf{y}_1, \mathbf{x}) - \mathbf{q}_{\theta}(\mathbf{y}_2, \mathbf{x})$ . Expanding the squared term in the objective of (9), the optimization problem reduces to

$$\min_{\theta} \quad \mathbb{E}_{\substack{\mathbf{x}, \\ \mathbf{y}_1, \mathbf{y}_2 \sim \mu}} \left[ (\mathbf{q}_{\theta}(\mathbf{y}_1, \mathbf{x}) - \mathbf{q}_{\theta}(\mathbf{y}_2, \mathbf{x}))^2 - 2[\mathbf{q}_{\theta}(\mathbf{y}_1, \mathbf{x}) - \mathbf{q}_{\theta}(\mathbf{y}_2, \mathbf{x})] \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right] \quad (13)$$

The issue lays in the unknown factor in the second term. Now comes to the IPO's peculiar derivations. The second term can be written as

$$\begin{aligned}
& \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2 \sim \mu} \left[ [q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})] \mathbb{E}_{\mathbf{y} \sim \mu} [v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \right] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x})) (v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}))] \\
&= \underbrace{\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})]}_{(*)} \\
&\quad - \underbrace{\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})]}_{(**)}
\end{aligned}$$

For (\*), notice that the two terms are essentially the same under expectation when  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}$  are i.i.d.

$$\begin{aligned}
(*) &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_2, \mathbf{y}} [q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_1, \mathbf{y}} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] + \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \\
&= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \tag{14}
\end{aligned}$$

The particular names of variables under expectation do not matter as long as they share the same distribution. A simple example demonstrating this point:

$$\mathbb{E}_{x \sim p} [f(x)] + \mathbb{E}_{y \sim p} [f(y)] = \mathbb{E}_{x \sim p} [f(x)] + \mathbb{E}_{x \sim p} [f(x)] = 2 \mathbb{E}_{x \sim p} [f(x)].$$

Similar trick can be applied for (\*\*) as well:

$$\begin{aligned}
(**) &= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x}) + q_\theta(\mathbf{y}_2, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}, \mathbf{x})] \\
&= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y} \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= 2 \mathbb{E}_{\mathbf{y}_1 \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})] \mathbb{E}_{\mathbf{y}_2, \mathbf{y} \sim \mu} [v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1 \sim \mu} [q_\theta(\mathbf{y}_1, \mathbf{x})] \tag{15}
\end{aligned}$$

where the last equality holds because  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}$  are independent, and  $\mathbb{E}_{\mathbf{y}_2, \mathbf{y} \sim \mu} [v(\mathbf{y}_2, \mathbf{y}, \mathbf{x})] = 0.5$  which in turn can be derived from conditions (8b) and (8c). Combining (14), (15) gives

$$\begin{aligned}
(*) - (**) &= 2 \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] - \mathbb{E}_{\mathbf{y}_1} [q_\theta(\mathbf{y}_1, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [2q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - (v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) + v(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x}))q_\theta(\mathbf{y}_1, \mathbf{x})] \quad (\text{by (8c)}) \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_2, \mathbf{y}_1, \mathbf{x})q_\theta(\mathbf{y}_1, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [q_\theta(\mathbf{y}_1, \mathbf{x})v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})q_\theta(\mathbf{y}_2, \mathbf{x})] \\
&= \mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})]
\end{aligned}$$

The expression in (13) then becomes

$$\begin{aligned}
& \arg \min_{\theta} \mathbb{E}_{\substack{\mathbf{x}, \\ \mathbf{y}_1, \mathbf{y}_2 \sim \mu}} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))^2] - 2 \mathbb{E}_{\substack{\mathbf{x}, \\ \mathbf{y}_1, \mathbf{y}_2 \sim \mu}} [(q_\theta(\mathbf{y}_1, \mathbf{x}) - q_\theta(\mathbf{y}_2, \mathbf{x}))v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x})] \\
&= \arg \min_{\theta} \mathbb{E}_{\substack{\mathbf{x}, \\ \mathbf{y}_1, \mathbf{y}_2 \sim \mu}} [(h_\theta(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) - v(\mathbf{y}_1, \mathbf{y}_2, \mathbf{x}))^2]
\end{aligned}$$



**“Notation”.** I am aware that I have used some terminologies quite freely and that could confuse rigorous readers. Here I like to clarify some of them

- By “learning”, I mean finding a mapping that can produce the right output given unseen input.
- By “policy”, I mean a language model
- By “generation”, I mean certain mechanism to sample  $\mathbf{y} \sim \pi(\mathbf{y} \mid \mathbf{x})$